

**Московский государственный технический
университет им. Н. Э. Баумана**

Курс «Технологии машинного обучения»

Отчёт по рубежному контролю №1

«Технологии разведочного анализа и обработки данных.»

Вариант № 5

Выполнила:
Буйдина К.А.
группа ИУ5-63Б

Проверил:
Гапанюк Ю.Е.

Дата: 11.04.25

Дата:

Подпись:

Подпись:

Москва, 2025 г.

Задание:

Номер варианта: **5**

Номер задачи: **1**

Номер набора данных, указанного в задаче: **5**

<https://www.kaggle.com/mohansacharya/graduate-admissions> (файл Admission_Predict.csv)

Для студентов групп ИУ5-63Б, ИУ5Ц-83Б - для произвольной колонки данных построить график "Ящик с усами (boxplot)".

Задача №1.

Для заданного набора данных проведите корреляционный анализ. В случае наличия пропусков в данных удалите строки или колонки, содержащие пропуски. Сделайте выводы о возможности построения моделей машинного обучения и о возможном вкладе признаков в модель.

Ход выполнения:

Задача 1

Для заданного набора данных проведите корреляционный анализ.

В случае наличия пропусков в данных удалите строки или колонки, содержащие пропуски.

Сделайте выводы о возможности построения моделей машинного обучения и о возможном вкладе признаков в модель.

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 df = pd.read_csv('Admission_Predict.csv')
6 df.drop('Serial No.', axis=1, inplace=True)
7 print("Количество пропущенных значений в каждой колонке:")
8 print(df.isnull().sum())
```

```
Количество пропущенных значений в каждой колонке:
GRE Score      0
TOEFL Score    0
University Rating  0
SOP            0
LOR            0
CGPA           0
Research       0
Chance of Admit  0
dtype: int64
```

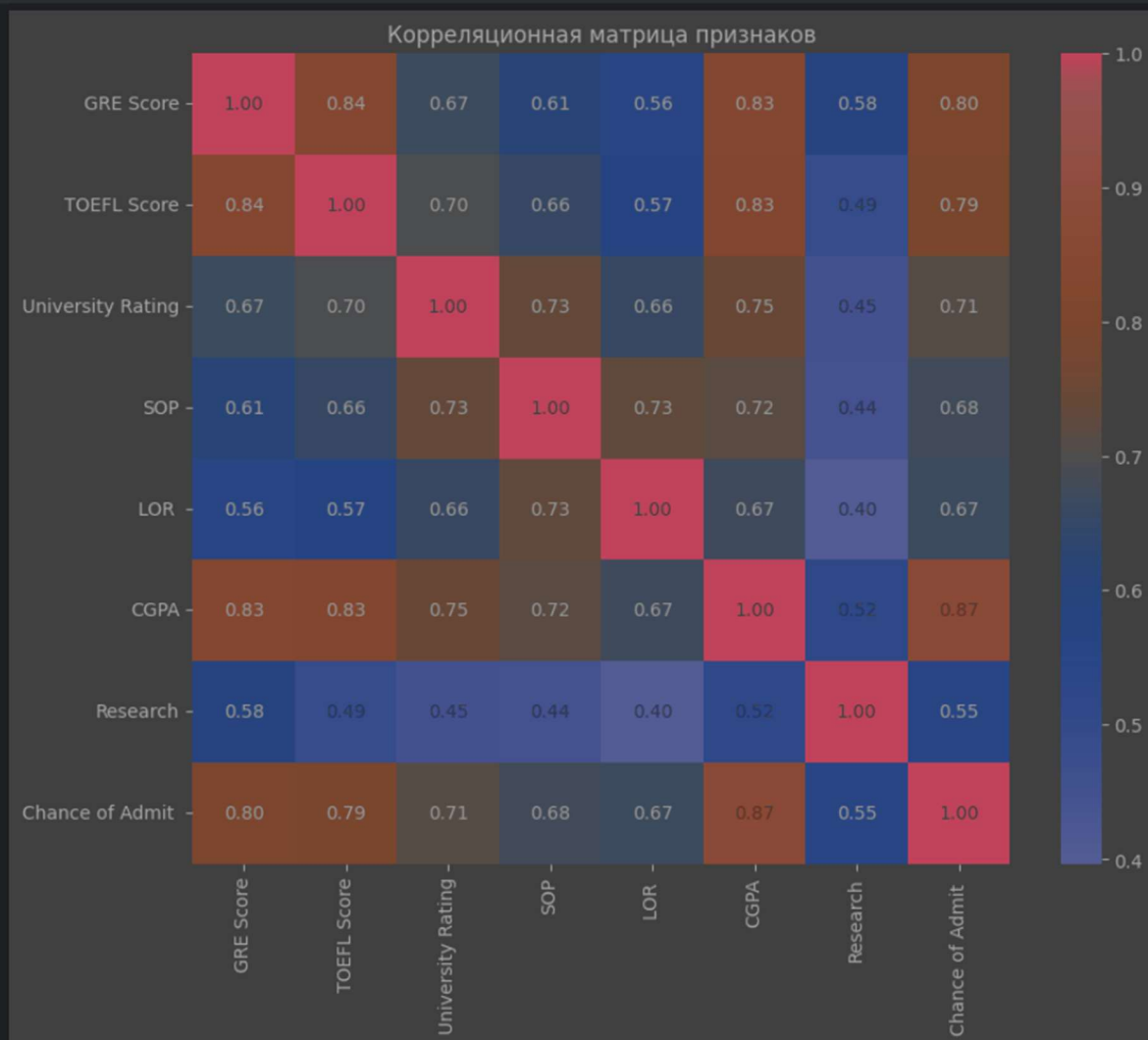
вывод - пропусков не было

```
1
2 # корреляция
3 correlation_matrix = df.corr()
4 plt.figure(figsize=(10, 8))
5 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
6 plt.title('Корреляционная матрица признаков')
7 plt.show()
8
```

```

1 |
2 # корреляция
3 correlation_matrix = df.corr()
4 plt.figure(figsize=(10, 8))
5 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
6 plt.title('Корреляционная матрица признаков')
7 plt.show()
8

```



Выводы о возможности построения моделей машинного обучения:

На основе корреляционной матрицы можно сделать следующие выводы:

- Целевая переменная 'Chance of Admit' имеет сильную положительную корреляцию с признаками 'GRE Score', 'TOEFL Score', 'University Rating', 'SOP', 'LOR' и 'CGPA'. Это говорит о том, что увеличение значений этих признаков приводит к увеличению вероятности поступления.
- Признак 'Research' также имеет положительную корреляцию с 'Chance of Admit', хотя и не такую сильную, как предыдущие.
- Между некоторыми входными признаками также наблюдается довольно высокая корреляция (например, между 'GRE Score' и 'TOEFL Score', между 'SOP' и 'LOR', а также между 'CGPA' и 'GRE Score'/'TOEFL Score'). Это называется мультиколлинеарность (корреляция между входными признаками), которую следует учитывать при построении некоторых моделей (например, линейной регрессии). Для моделей, основанных на деревьях решений (например, Random Forest, Gradient Boosting), мультиколлинеарность обычно не является серьезной проблемой, то есть она в целом не принципиальна и если она есть, то хуже не будет.

Возможный вклад признаков в модель:

- Наибольший вклад в модель, вероятно, внесут признаки с самой высокой корреляцией с 'Chance of Admit': 'CGPA', 'GRE Score', 'TOEFL Score', 'SOP' и 'LOR'.
- 'University Rating'
- Признак 'Research' - наименьший вклад

```
1 |
2 # boxplot
3 plt.figure(figsize=(8, 6))
4 sns.boxplot(y=df['CGPA'])
5 plt.title('Ящик с усами для CGPA')
6 plt.ylabel('CGPA')
7 plt.grid(True)
8 plt.show()
```

